

Docket Number: POU920020016US1

METHOD FOR DETECTING THE QUICK  
RESTART OF LIVENESS DAEMONS IN A  
DISTRIBUTED MULTINODE DATA  
PROCESSING SYSTEM

APPLICATION FOR  
UNITED STATES LETTERS PATENT

"Express Mail" Mailing Label No.: ET089969005US

Date of Deposit: February 15, 2002

I hereby certify that this paper is being  
deposited with the United States Postal Service as  
"Express Mail Post Office to Addressee" service  
under 37 CFR 1.10 on the date indicated above and is  
addressed to Box Patent Application, Assistant  
Commissioner for Patents, Washington, D.C. 20231.

Name: Susan L. Phelps

Signature: Susan L. Phelps

INTERNATIONAL BUSINESS MACHINES CORPORATION

**Method for Detecting the Quick Restart of Liveness Daemons  
in a Distributed Multinode Data Processing System**

**Background of the Invention**

The present invention is generally directed to insuring the continuation of consistent group formation events in a distributed topology liveness system, that is, in a multinode data processing system in which node and/or adapter liveness is communicated throughout the system via heartbeat messages, which are messages that are sent periodically and which indicate node and/or adapter liveness. More particularly, the present invention is directed to a method for detecting a situation in which a liveness daemon running on one of the nodes has been subject to a rapid restart. Even more particularly, the present invention is directed to a method for determining the existence of such quick restart events and for providing a proper indication thereof to other nodes within the network, with the particular objective of avoiding grouping inconsistencies which are situations in which one node set sees another node set fail in some way without the other node set being aware of the fact that the first node set has also failed. In short, all of the nodes within a node set should have the same view as to the operating status of the other nodes in the node set.

A proper understanding of the present invention is best obtained from an appreciation of the environment in which it is intended to operate. The present invention is employed in multinode data processing systems. These systems include a plurality of nodes each of which incorporates a data processing element which is coupled locally to its own memory system which typically includes both a volatile random access memory and a nonvolatile random access memory. The volatile memory typically comprises an array of semiconductor memory chips. The nonvolatile memory typically comprises a rotating magnetic or optical storage device. The data processing element also typically comprises a central processing unit (CPU). Each node includes one or more data processing elements. The nodes also include adapters which are communications devices which permit messages to be sent from one node to another node or to a plurality of other nodes. Internodal communications typically take place through a switch device which routes transmitted messages to destination nodes within the system.

In order to carry out various data processing functions, the nodes within any given multinode network are organizable into sets of nodes. Nodes and/or their associated adapters sometimes experience problems, delays or failures. Accordingly, from time to time during the operation of individual nodes, system checks are undertaken to make sure that the nodes are still alive and functioning. This checking is performed via heartbeat message transmissions. Each node in the system is assigned one or more "downstream" nodes for the purpose of periodically sending a message indicating liveness status. In preferred embodiments, heartbeat signals are only sent to a single other node. However, it is quite easy to instead employ a predefined list of node destinations for receipt of heartbeat signals from any or all of the nodes in the network. These liveness message transmissions are handled by daemon programs running on the various nodes in the system.

10  
POU920020016US1  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25

Distributed multinode data processing systems of the kind contemplated herein employ heartbeat messaging protocols which are used to control group membership which, of course, shifts over time. It is control of the membership process to which the present invention is directed. This membership process typically includes the establishment of one of the nodes in a group as the so-called Group Leader (GL). The Group Leader acts as a coordinator for nodes coming into (joining) or for nodes exiting the group. Additionally, in the event that there is a problem with the Group Leader, there is preferably also a designated second node which is intended to act as a replacement for the Group Leader in the event that the Group Leader experiences a failure. This second, backup Group Leader is referred to as the Crown Prince. In the context of the present invention, the Group Leader and Crown Prince are employed in the "liveness" (heartbeating) layer. The present invention should not be confused with group membership services which are provided to "end user applications." In accordance with the present invention, "group membership," as referred to above, refers to the list of members in an Adapter Membership Group which occurs on each network being monitored. On the other hand, "node reachability" refers to the set of nodes that are considered to be alive, taking all of the adapter membership groups into consideration. In particular, it is noted that the notion of "node reachability" may include message hops through indirect paths that may cross network

boundaries. This set of nodes is supplied from the "liveness layer" to the "group communications layer" which runs on top of the "liveness" layer.

5

More particularly, the present application is concerned with two different scenarios which present potential problems with respect to group membership consistency across the nodes of the system or network. Accordingly, there is provided a method for determination of adapter and node death in a distributed system whereby node events are made consistent, that is, when a first node sees another node as being "down," the second other node, if alive, is still able to see the first node as being "down" within finite amount of time. When a node actually suffers a "permanent" crash the heartbeat mechanism, together with the associated "join" protocol, is able to provide sufficient control and communications amongst the remaining nodes to assure maximum functionality. Accordingly, the present invention does not come into play when nodes crash, since the basic heartbeat mechanism is able to cope with this situation; nonetheless, the present invention becomes important when communication failures and process blockages result in temporary loss of contact amongst a set of distributed peers in the liveness determination subsystem. The present method addresses two possible scenarios which could lead to inconsistent node grouping situations: (1) a node where the liveness daemon is stopped and restarted quickly; and (2) a node whose communications with the rest of the nodes suffers a temporary interruption.

20

25

In situations in which the liveness daemon running on one of the nodes is stopped and restarted in a short period of time, certain consistency problems can be engendered. For example, typically it happens that when the liveness daemon restarts, for each local adapter, a message is transmitted which "proclaims" the existence and the willingness of the sending node to become a group leader; it is, in generic terms, a request to know which other nodes are "out there." These aspects are discussed in more detail below where the nature of the "PROCLAIM" message is considered. However, the other nodes in the group still consider the restarting node (and/or adapter) as being part of the previous group. Accordingly, group membership is no longer consistent in the sense that there is a lack of symmetry among the various nodes with regards to the "known" status of the other nodes. When this situation is caused by the "quick"

restart of the liveness daemon, it is referred to herein as the "bouncing node" problem or scenario.

Likewise, a problem can occur if a first node, say Node 1, has a temporary communication problem. If the problem lasts long enough for the other nodes to expel Node 1 from the group, but not long enough for the local adapter to be declared down, the other nodes can form a new Adapter Membership Group, G2, while the adapter at Node 1 is still considered as being part of the previous group, G1 (which contains all the adapters). The adapter at Node 1 then attempts to dissolve the group, since it will have gotten no answer to a liveness ("DEATH") message that it sent when its old upstream neighbor stopped sending heartbeat signals to it. (For a discussion of a more specific and preferred characterization of the notion of dissolving a group, attention is directed below to Section 2.2). Upon "dissolving" the group, the adapter at Node 1 reinitializes into a "group" with only a single node, which is referred to herein as a singleton group and it resumes operation. Singleton groups are inherently unstable groups since they are typically destined to soon experience a change to inclusion in a larger group. If this all happens before the adapter on Node 1 is able to form a stable group, then Node 1 never sees any "node down" events, where the other nodes see Node 1 as being "down," especially if this is the only adapter group to which Node 1 belongs. Accordingly, the recognition of this problem brings along with it the notion that some groups are more stable (from time to time) than other groups, and that special handling is required to insure group membership consistency across the network.

20

### Summary of the Invention

In accordance with a preferred embodiment of a first aspect of the present invention, there is provided a method for detecting the quick restart of liveness daemons in a distributed, multinode data processing system in which the nodes communicate liveness indicia in the form of heartbeat messages via adapters coupled to each node. In this method a first message (PROCLAIM) is sent from a first node to other nodes in the network that do not yet belong to the local node's adapter membership group. This message contains some indicia that the sending node has recently experienced an adapter restart. This information, together with locally stored group membership

information, is used to determine that a quick restart has actually occurred at the sending node. This situation is handled by expelling the node from the group as a means for insuring correct group membership status.

5 In accordance with a preferred embodiment of a second aspect of the present invention, there is provided a method for detecting node reachability inconsistencies in the presence of temporary node communication failures or temporary daemon blockage. To accomplish this, an indication of a last stable adapter membership group is maintained at each node. The group join protocol is thus enabled to provide a PREPARE\_TO\_COMMIT (PTC) message which includes a flag which indicates that the message recipient is considered as belonging to the same stable membership group as the message sender. As used herein, the term "stable" refers to a characteristic for a node or node group which implies that there is only a small likelihood that group membership for that node or node group will change in the near future. In particular, nodes that find themselves isolated as the only members of a group try to join a group as soon as possible. Therefore, such singleton groups are referred to herein as being inherently unstable.

10 Accordingly, it is an object of the present invention to insure the existence of consistent group membership information across a plurality of nodes in a distributed, multinode data processing system.

15 It is also an object of the present invention to provide a mechanism to guard against group membership inconsistencies which might arise as the result of the failure and quick restart of a node and/or one of more of its associated adapters.

20 It is a further object of the present invention to employ existing group membership control protocols as a mechanism for communicating proper group membership status.

It is yet another object of the present invention to insure proper group membership status in the face of temporary node communication failures.

It is a still further object of the present invention to insure proper group membership status in the face of temporary node daemon blockages.

It is also an object of the present invention to provide an indication that a node is included within a stable group.

5 It is yet another object of the present invention to provide an indication of stability for the nodes in a data processing network.

It is also an object of the present invention to increase the reliability and availability of distributed data processing systems.

10 It is also an object of the present invention to expand the capabilities of Topology Services in terms of its utility without significantly altering any of its application programming interfaces (APIs) or its protocols.

15 Lastly, but not limited hereto, it is an object of the present invention to provide a cooperative relation between the first and second aspects of the present invention (quick restart and failed communications, respectively) to particularly address the problem of assuring consistent node viewpoints with respect to adapter group membership and node reachability.

20 The recitation herein of a list of desirable objects which are met by various embodiments of the present invention is not meant to imply or suggest that any or all of these objects are present as essential features, either individually or collectively, in the most general embodiment of the present invention or in any of its more specific embodiments.

#### Brief Description of the Drawings

The subject matter which is regarded as the invention is particularly pointed out and distinctly claimed in the concluding portion of the specification. The invention, however, both as to

organization and method of practice, together with the further objects and advantages thereof, may best be understood by reference to the following description taken in connection with the accompanying drawings in which:

5      Figure 1 is a schematic block diagram illustrating the connections of multiple nodes to one or more networks through a plurality of adapters;

Figure 2A is a schematic block diagram illustrating the transmission of a PROCLAIM message as part of a join protocol;

Figure 2B is a schematic block diagram illustrating the JOIN response to a PROCLAIM message;

Figure 2C is a schematic block diagram illustrating the transmission of the PREPARE\_TO\_COMMIT (PTC) message as part of a group joining protocol;

Figure 2D is a schematic block diagram illustrating the transmission of the PTC\_ACK (prepare to commit acknowledgment) message in response to the PTC message;

15      Figure 2E is a schematic block diagram illustrating the transmission of the COMMIT\_BCAST (broadcasted commit) message;

Figure 2F is a schematic block diagram illustrating the transmission of an acknowledgment to the COMMIT\_BCAST message;

20      Figure 2G is a schematic block diagram illustrating the new group formed as a result of the messages transmitted in Figures 2A through 2F and further illustrates the flow graph for the subsequent transmission of heartbeat messages;

Figure 3A is a schematic block diagram illustrating the heartbeat message path at the beginning of the protocol dealing with node and/or adapter death (failure);

Figure 3B is a schematic block diagram illustrating the transmission of a DEATH message to the Group Leader in the event that the node at the lower right fails to pass along its heartbeat message;

Figure 3C is a schematic block diagram illustrating the beginning of the chain of events which follows the receipt of a DEATH message, namely the transmission of the PREPARE\_TO\_COMMIT message;

Figure 4A is a schematic block diagram illustrating transmission to the Group Leader of the NODE\_CONNECTIVITY message;

Figure 4B is a schematic block diagram illustrating transmission to the Group Leader of the GROUP\_CONNECTIVITY message;

Figure 4C is a schematic block diagram illustrating group connectivity across two networks and the forwarding of group connectivity information;

Figure 5A is a schematic block diagram illustrating an initial state for the Node Connectivity Table just prior to the occurrence of a node failure (death);

Figure 5B is a schematic block diagram illustrating the death of Node 2 and the formation of Nodes 1, 3 and 4 into Adapter Membership Group A\_2;

Figure 5C is a schematic block diagram illustrating the propagation of the Group Connectivity Message to all nodes;

Figure 6 is a time flow diagram which compares the activities in different nodes so as to illustrate consistency problems arising as a result of the quick restart of a daemon on Node 1;

Figure 7 is a time flow diagram which compares the activities in different nodes so as to illustrate consistency problems arising from temporary adapter communication failure;

5       Figure 8 is a time flow diagram which compares the activities in different nodes so as to illustrate consistency problems arising from a temporary adapter communication problem in which the problem node is the Group Leader or Crown Prince;

Figure 9 is a schematic block diagram illustrating the structure of adapter IDs and Group IDs;

Figure 10 is a block diagram illustrating the preferred format for the protocol message packets that are sent over the network;

Figure 11 is a schematic block diagram illustrating a sample structure for the adapter and group IDs when the daemon at Node 1 terminates and is restarted;

Figure 12 is a schematic block diagram illustrating the processing that occurs when a live node detects that a remote daemon has restarted;

15       Figure 13A is a schematic block diagram illustrates message flow that occurs when a daemon that is restarted detects the fact that a previous instance used to belong to an Adapter Membership Group because of the heartbeat messages that it receives while it is in a singleton group;

Figure 13B is the second portion of Figure 13A;

Figures 14A through 14D are a sequence of block diagrams illustrating a solution to the problem of communication interruption that occurs when Node 3's adapter suffers a temporary failure;

5       Figure 14B is a schematic block diagram following Figure 14A in time and more particularly illustrating the point in time that Node 2 commits to a new Adapter Membership Group while Node 3 is still in the process of missing heartbeat messages from its neighbor;

Figure 14C is a schematic block diagram following Figures 14A and 14B in time and more particularly illustrating the point in time that Node 3 sends a PTC message when it stops receiving heartbeat messages from its upstream neighbor and also particularly showing the rejection of the PTC messages because of discrepancies in the last\_stable\_group results;

Figure 14D is a schematic block diagram following Figures 14A, 14B and 14C in time and more particularly illustrating the fact that Node 3 is forced into a singleton group at which point it updates its last\_stable\_group indicator;

15       Figures 15A through 15E are a sequence of block diagrams similar to Figures 14A through 14D which particularly illustrate the situation in which Nodes 1, 2 and 3 are part of the same Adapter Membership Group;

20       Figure 15B is a block diagram following Figure 15A in time and more particularly illustrating the situation in which, as a result of a temporary failure in Node 1, Node 3 commits to a new Adapter Membership Group while Node 1 is still in the process of missing heartbeat messages from its upstream neighbor;

Figure 15C is a block diagram following Figure 15B in time which illustrates the course of action, following that shown in Figure 115B, in which Node 1 dissolves its group and forms a singleton (and thus an unstable) group;

Figure 15D is a block diagram following Figure 15D in time which illustrates the course of action, following that shown in Figure 15C, in which Node 3 sends a PTC message when Node 1 responds to the PROCLAIM message with a JOIN message; and

5       Figure 15E is a block diagram illustrating the course of action, following that shown in Figure 15D, in which since Node 3 does not get replies to its PTC messages, it is eventually forced to form a singleton group.

### Detailed Description of the Invention

#### **1.0 Controlling Group Membership**

Adapter and node liveness determination lies at the heart of any highly available distributed cluster data processing system. In order to provide high availability services, a cluster system should be able to determine which nodes, networks, and network adapters in the system are working and be able to accurately determine the group or groups to which they belong. The failure in any such component should be detected as soon as possible and indications of such failure should be passed along to higher level software subsystems for recovery processing by the cluster software and/or applications running on the cluster.

20       Determination of node, network, and network adapter liveness is often made through the use of daemon processes running on each node of the distributed system. Daemons run distributed protocols and exchange liveness messages that are forced through different network paths in the data processing system. If no such liveness messages are received within a predetermined time interval, then the sending node or network adapter is assumed to be not working ("dead") by the other nodes.

Any method of liveness determination can be subjected to "false down" events, where nodes or network adapters are incorrectly notified as being down or unreachable. Such false events may happen, for example, when temporary communication failures prevent the liveness messages

from reaching their destination(s). False "node down" events may also happen when the liveness determination daemon is prevented from being scheduled because of CPU scheduling, memory contention, excessive interrupts, and other factors. The daemon being stopped is yet another source of false "node down" notifications.

5 In the presence of these false events, it is important to provide consistent node reachability notifications: when a node sees the other node as down, the other node -- if alive -- should see the first as down within a finite, preferably predetermined, time interval. The absence of such consistency may lead to undesirable effects, since software layers above the liveness determination "layer" may be unable to reach an agreement regarding the topology's health (that is, the configuration of nodes within a group together with identifiable paths by which one node may be reached from another), with different nodes having different views of which nodes are reachable.

## 2.0 Heartbeat Protocols

20 To explain the mechanisms of the present invention, and how they are employed in Topology Services (a set of system utility programs and defined API calling structures), adapter membership ("heartbeating") protocols in the subsystem are explained herein in some detail. Topology Services is the layer in the infrastructure which is responsible for detecting the health of adapters, nodes, and networks.

In order to monitor the health and connectivity of the adapters in each network, all adapters in  
25 the network attempt to form at least one "Adapter Membership Group" (AMG), which is a group containing all network adapters within the network that can communicate with each other. Adapters in an AMG monitor the "liveness" of each other. When an AMG is formed, all group members receive an "AMG id" (that is, a unique group membership identifier) which identifies the AMG. If, at some point in time an adapter fails, it is expelled from the group, and new adapters that are powered up are invited to join the group. In both cases, a new AMG with a new "AMG id" is formed. Each AMG has one member that is the Group Leader (GL), and all

members know who the Group Leader is. Note that a node may belong to several AMGs, one for each of its (network) adapters.

Each AMG has an id, which is included in all protocol messages. The group id includes the GL identification (chosen to be its Internet Protocol (IP) address) and an instance number 5 (chosen to be the time stamp of which indicates when the AMG was formed). Note that the group id is chosen to be the IP address for convenience and that any scheme for assigning a unique and sequentially orderable identifier may be employed.

Each member of an AMG also has an id, which includes the member identification (chosen to be its IP address) and an instance number (chosen to be the time stamp of when its daemon was started or when its adapter was reinitialized).

To determine the set of adapters that are alive in each network, an adapter membership protocol is run in each of the networks. Messages in this protocol are sent using UDP/IP ("User Datagram Protocol" / "Internet Protocol").

Adapters that are alive form an AMG, where members are organized in a virtual ring topology. To ensure that all group members are alive, each member periodically sends "HEART 15 BEAT" messages to its "downstream neighbor" and monitors "HEART BEAT" messages from its "upstream neighbor." Protocols are run when adapters fail or when new adapters become functional. The goal of such protocols is to guarantee that the membership group contains at each moment all (and only) the adapters in the network (but only those belonging to the cluster) 20 that can communicate with each other.

Besides the Group Leader, each group has a "Crown Prince" (backup group leader). The group leader is responsible for coordinating the group protocols, and the Crown Prince is responsible for taking over group leadership if the group leader, or its adapter, fails. Both the choice of Group Leader and Crown Prince, and the position of the adapters in the ring, are 25 determined by a predefined adapter priority rule, which is typically chosen to be the adapters' IP

address, hence the desire, as stated above, that its indicia be able to provide a sort into a unique ordering sequence.

A list of all possible adapters in each network is contained in a configuration file that is read by all of the nodes at startup and at reconfiguration time.

5

## 2.1 Join Protocol

In order to attract new members to the group, the Group Leader in each group periodically sends "PROCLAIM" messages to adapters that are in the adapter configuration but do not currently belong to the group. The message is only sent to adapters having a lower IP address than that of the sender. See Figures 2A through 2G and the discussion in Section 7 below.

The "PROCLAIM" messages are ignored by all adapters that are not group leaders. A Group Leader node receiving a "PROCLAIM" message from a higher priority (higher IP address) node responds with a "JOIN" message on behalf of its group. The message contains the membership list of the "joining group."

15 A node GL1 (Group Leader #1) receiving a "JOIN" message from GL2 (Group Leader #2) attempts to form a new group containing the previous members plus all members in the joining group. GL1 then sends a "PTC" ("Prepare To Commit") message to all members of the new group, including GL2.

20 Nodes receiving a "PTC" message reply with a "PTC\_ACK" message. All nodes from which a "PTC\_ACK" message is received are included in the new group. The group leader (GL1) sends a "COMMIT" message, which contains the entire group membership list, to all new group members.

Receiving a "COMMIT" message marks the transition to the new group, which now contains the old members plus the joining members. After receiving this message, a group member starts

sending "HEART BEAT" messages to its (possibly new) downstream neighbor, and starts monitoring "HEART BEAT" messages from its (possibly new) upstream neighbor.

Both "PTC" and "COMMIT" messages require an acknowledgment to ensure they have been received. If no acknowledgment is received then a finite number of retries is made. Failure to respond to a "PTC" message -- after all retries have been exhausted -- results in the corresponding adapter not being included in the new group. If a liveness daemon fails to receive a "COMMIT" message after all retries of the "PTC\_ACK" message, then the local adapter gives up the formation of the new group and reinitializes itself into a singleton group. This phenomenon should only occur in the relatively rare case where the Group Leader fails in the short time window between sending the "PTC" message and the "COMMIT" message.

When the Topology Services daemon is initialized, it forms a singleton adapter group (of which the node is the Group Leader) in each of its adapters. The node then starts sending and receiving "PROCLAIM" messages.

## 2.2 Death Protocol

A node or adapter monitors "HEART BEAT" messages coming from its "upstream neighbor" (the adapter in the group that has the next highest IP address among the group members). When no "HEART BEAT" messages are received for some predefined period of time, the "upstream neighbor" is assumed to have failed. A "DEATH" message is then sent to the group leader, requesting that a new group be formed. See Figures 3A through 3C and the discussion in Section 7 below.

Upon receiving a "DEATH" message, the group leader attempts to form a new group containing all adapters in the current group except the adapter that was detected as failed. The group leader sends a "PTC" message to all members of the new group. The protocol then follows the same sequence as that described above for the Join protocol.

After sending a "DEATH" message, the daemon expects to receive a "PTC" message shortly. A number of retries is attempted, but if no "PTC" message is received then the interpretation is that the GL adapter (or its hosting node) died and that the "Crown Prince" adapter also died, and therefore was unable to take over the group leadership. In this case the adapter reinitializes itself into a singleton group and also sends a "DISSOLVE" message, inviting all group members to do the same. This is the mechanism that allows all members of the group to find out about the simultaneous demise of the Group Leader and Crown Prince member nodes.

### 2.3 Node Reachability

A node reachability protocol is used to allow computation of the set of nodes that are reachable from the local node (and therefore considered alive). Since not all nodes may be connected to the same network, some nodes may be reachable only through a sequence of multiple network hops. Complete node reachability determinations can only be computed when information about all networks, even those that do not span all nodes, is taken into account.

To compute node reachability, an eventual agreement protocol is used: reachability information at each network is propagated to all networks; when the network topology stops changing, eventually all nodes have consistent information about all networks. Each node is then be able to compute the set of reachable nodes independently and arrive at a consistent result.

Periodically, and until some stopping criteria instruct the daemon to stop doing so, the nodes send the following messages:

a "Node Connectivity Message" (NCM) which is sent from all group members to the Group Leader. A Node Connectivity Message for a given network contains the AMG id for that network plus all of the "disabled AMG ids" for the local adapters that are disabled. A node sends NCMs to each Group Leader of the groups to which the local adapters belong. The Group Leader stores all the information coming from the NCM's in a "Node Connectivity Table" (NCT). The

NCT stores the (local view of the) global network topology, and contains the AMG id for each node and network adapter in the system. Any two nodes that have the same AMG id are assumed to be connected to each other by the same network; and

5 a "Group Connectivity Message" (GCM) which is sent from each Group Leader to all group members. The GCM contains the AMG id and the list of nodes that belong to the AMG. Also, for each of these nodes, a list of all "disabled AMG ids" (in the other networks) is included. The information needed to send the GCM is extracted from the Group Leader's Node Connectivity Table. A node that receives a GCM updates its own Node Connectivity Table with the information in the message. If a daemon receiving a GCM notices that there are some groups to which the local adapters belong, whose members have not received that GCM, the daemon forwards the GCM to these groups. The goal is to propagate the GCM to all of the nodes in the system, even those that are not directly connected to the network that originated the GCM. Notice that the information sent in an NCM and GCM is a subset of the sender's NCT.

See Figures 4 and 5, and the discussion in Section 7 below.

#### **2.4 - Stable/Unstable AMGs**

To prevent "panic" actions of the protocol -- such as those caused by the absence of a  
20 "COMMIT" after all "PTC ACKs" or by the simultaneous failure of the Group Leader and the Crown Prince -- from causing major node reachability ripples, the concept of "stable" and "unstable" AMGs is now defined. Stable AMGs are those where steady state operations are occurring, while unstable AMGs are those where membership changes are still likely to occur (such as for singleton groups).

At initialization, singleton Adapter Membership Groups start in the unstable state, since it is expected that the adapter will join other peers in bigger groups. The change into a stable group occurs after an inactivity period where membership changes stop occurring. Once a group is stable, it remains stable until the adapter is forced to reinitialize itself because of a "panic" action.

AMG stability is tied to the Node Connectivity Table and to sending the Node Connectivity Message and the Group Connectivity Message in the following way: to prevent unnecessary node reachability ripples, no information about unstable groups is sent in NCMs and GCMs. This effectively removes unstable groups from the computation of node reachability, and has the desirable effect of eliminating the knowledge of some temporary membership changes from the software layers above.

### **3.0 Node Event Inconsistency**

#### *3.1 - Inconsistency Caused by Quick Restart*

A liveness daemon which stops (on request or due to a problem) and is then quickly restarted offers a chance for inconsistency to occur. In order to best appreciate this phenomenon, consider the following sequence of events which occurs when a daemon is stopped and then quickly restarted:

(1) the daemon at Node A is stopped;

(2) the daemon at Node A restarts and, for each local adapter, initiates sending "PROCLAIM" messages;

(3) the other nodes still consider Node A as part of the previous group. This situation continues until Node A is detected as dead on each of the AMGs due to lack of "HEART BEAT" messages coming from Node A;

- (4) Node A is finally detected as dead in each AMG and is expelled from the group; and
- (5) Node A is then finally allowed to rejoin each AMG.

The delay in "(3)" causes the following problems:

(1) reintegration of the "bouncing" node is seen as occurring too slowly; and

5 (2) if different networks have very different detection times, it is possible that Node A may be detected as being down and thereafter rejoins one of the groups before being ever detected as down in another network (which has a longer detection time). The net result is that, when node reachability is computed by the other nodes, Node A is never seen as going down at all.

10 The problem with the scenario in (2) above is that the daemon that restarted starts anew, with no memory of the previous AMG. If other nodes never detect that the node "failed," then they cannot take actions to integrate the node into the higher level node group.

### *3.2 - Inconsistency Caused by Quick Communication Interruption*

15 Some node event inconsistency problems are possible because of the inherent behavior of the base adapter membership protocols. The following are two examples of scenarios that could lead to inconsistent events.

20 (1) Node 1 has a temporary problem in its adapter. The problems lasts long enough for the other nodes to expel Node 1 from the group, but not long enough for the local adapter to be declared down. While the other nodes form a new AMG G2, the adapter at Node 1 initially considers itself still as part of the previous G1 (which is assumed in this example to contain all of the adapters). The adapter at Node 1 then attempts to dissolve the group, since it got no answer to its "DEATH" message that it sent when its old upstream neighbor stopped sending heartbeat messages to it. Upon "dissolving" the group, the adapter at Node 1 then reinitializes into a

singleton unstable group and resumes operation. If the adapter is working again, "PROCLAIM" messages eventually arrive, and the adapter is brought back into the group. If this all happens before the adapter on Node 1 can form a stable group, then Node 1 never sees any node down events, whereas the other nodes will have seen node 1 as down if this is the only adapter group to which Node 1 belongs.

(2) This next example is similar to the one above, but this time it is assumed that Node 1 used to be the Group Leader. During the temporary outage, other adapters in the AMG form group G2 and expel the adapter at Node 1. Node 1 only perceives that it was expelled from the group when the heartbeats from its upstream neighbor stop coming. At some point, Node 1 declares the upstream neighbor dead and simply sends a "PTC" message to its old group membership. The other nodes, upon seeing the "PTC" message from an adapter with higher priority, immediately respond to the "PTC" message, and a new group G3 is formed. While the other nodes will have seen Node 1 failing and then coming back, Node 1 does not actually see the others failing (except possibly for its old upstream neighbor). Node 1 is completely oblivious to being expelled from the adapter group.

## 4.0 Detection of Bouncing Nodes

The detection of "bouncing nodes" (that is, nodes where the liveness daemon exits for any reason and is then restarted within a short period of time) is based on the bouncing nodes and the live nodes finding about the bounce by using normal liveness protocol messages.

20

### 4.1 Live Nodes Detect Bouncing Nodes

One way by which the current nodes in the group can detect bounced members is by receiving "PROCLAIM" messages from them. The "PROCLAIM" message can indeed reveal that the source of the message is a bounced entity by determining that all three of the conditions indicated below exist:

(1) the group id of the message is different from the recipient's group id (since the bouncing daemon is definitely now using a different group id -- at least the instance number part);

(2) the IP address of the sender of the message is still listed as part of the current group; and

(3) the instance number of the sender's adapter id is different from the one listed as group member (the instance number changes when the adapter is reinitialized or when the daemon starts).

If a "PROCLAIM" message is received where all three of the conditions listed above are true, then the assessment is that the message came from a group member that bounced. To speed up the detection of the bounce and to allow faster reintegration of the bouncer, the best course of action is to expel it from the group, which can be done by sending a "DEATH" message for the bouncing adapter.

Since the "PROCLAIM" message is likely to reach all group members, then all of them would try to send a "DEATH" message for the bounced adapter, which is wasteful. The alternative is for only the bouncer's downstream neighbor to send the message. Accordingly, such a process is indicated in the pseudo-code provided below:

#### Handling a "PROCLAIM" message

```
if( from_group != from my group id  &&
    from_IP_address is still part of my group  &&
    from_Instance != from the id that is in the group  &&
    I am the id's downstream neighbor ) {
    send a "DEATH" message for id
}
```

One additional method for detecting a bounced daemon includes a step wherein a bounced daemon sends a "JOIN" message even before the Group Leader is notified about the demise of the adapter.

Though "PROCLAIM" messages are usually sent only to adapters which are not currently part of the AMG, implementations of the present protocol may, if desired, use a hardware broadcast to send the message, in which case even members of the AMG may receive the message. In this case, the Group Leader receives a "JOIN" message from an adapter which is still member of the group. This situation can be recognized as a "bounce" by the GL, which then invokes the PTC-COMMIT sequence to expel the adapter from the group.

#### *4.2 - Bounced Nodes Detect That They Bounced*

Normally, a daemon that bounces starts anew with no memory of a previous instance. On the other hand, a bounced daemon that used to be part of a group is likely to receive "HEART BEAT" messages from its old "upstream neighbor." Such "HEART BEAT" messages tell the bouncing daemon that it bounced quicker than the time it takes to detect a remote adapter as dead.

Again, the goal is to cause the bouncing adapter to be expelled from the previous group as soon as possible. The first thought which occurs as a method for accomplishing this goal is for the daemon that receives such a "HEART BEAT" message to send a "DEATH" message for itself, but this does not work because the bouncing daemon does not know who the Group Leader is, and therefore does not know to whom to send the "DEATH" message. In addition, the Group Leader may have itself been the recipient of the message (that is, the node that bounced). The solution to this problem is for the bouncing daemon to send a new "NOT YOUR NEIGHBOR" message back to the sender of the "HEART BEAT" message. The recipient of this message, being part of the previous group and knowing who the Group Leader is, reacts by sending a "DEATH" message to the Group Leader. Accordingly, such a process is indicated in the pseudo-code provided below:

Receiving a "HEART BEAT" message which is not to the current group:

```
if ( I am part of a singleton group ) {  
    Reply with "NOT YOUR NEIGHBOR" message  
}
```

5 Receiving a "NOT YOUR NEIGHBOR" message:

```
if (sender's IP address is part of my group ) {  
    Find id corresponding to the IP address  
    Send a "DEATH" message to GL for the id  
}
```

## 5.0 Solution to the Quick Communication Interruption Problem

The steps described below are carried out to address the two situations cases described above in section 3.2. The object of performing these steps is to force both sides of a merging group to see roughly the same events prior to the merge.

(1) Each node keeps, for each local adapter, a copy of the last stable AMG to which the local adapter belongs ("last\_stable\_group"). The rationale for keeping only the stable groups is that only stable groups result in the desire for node reachability to be recomputed.

(2) When sending a "PTC" message, the sender adds an "in\_previous\_stable\_group" flag to the message, according to whether the destination belonged to the last\_stable\_group AMG -- usually the previously formed AMG prior to the new group being committed.

20 (3) When processing a "PTC" message, a node handles two pieces of information: (1) the "in\_previous\_stable\_group" in the message; and (2) whether the sender of the message belongs to the receiver's "last\_stable\_group" group. Unless these 2 pieces have the same TRUE/FALSE value, the "PTC" message is rejected.

The mechanism above withstands both examples in section 3.2 above and is also effective in more normal cases, such as when two AMGs merge in a PROCLAIM-JOIN-PTC sequence. See section 6.2 below.

## 6.0 - Sample Scenarios

5 Some scenarios are presented to depict how the protocols presented herein work to effect their various purposes. In the case of multiple bouncing two of the separate protocols described herein work together in a cooperative fashion to further assure consistent group membership.

### 6.1 - Detection of Bouncing Nodes

#### (a) Single bounce

After a daemon "bounces," either of the bouncing detection mechanisms should be activated, in any case resulting in the bouncing "adapter" being removed from the group. Since this happens in all AMGs more or less simultaneously, the node is effectively detected as dead by the others before it can rejoin the AMGs.

#### (b) Multiple bounces

15 This example considers the case wherein there are multiple bouncing adapters, that is, when a number of nodes bounce, while others fail. The concern here is what happens when a bouncing daemon quickly joins in a group with some other adapter whose daemon also bounced. Not being the Group Leader of the group, the daemon does not send any "PROCLAIM" messages. In addition, when multiple nodes bounce it may happen that a bouncing daemon's upstream 20 neighbor has also bounced. Therefore the usual mechanisms are not active in causing the bouncing adapter to be expelled.

This situation is salvaged by the methods of the present invention through the observation that at least one of the bouncing daemons becomes the Group Leader in its new group; the others

might become members of this very same group. The Group Leader sends "PROCLAIM" messages periodically, resulting in a "DEATH" message being sent for it. The Group Leader of the original group then attempts to form a new group, but then none of the bouncing daemons should reply to the "PTC" message. This happens because the quick communication interruption mechanism described in section 5.0 above comes into play: the bouncing daemon is still part of the Group Leader's group, while the Group Leader itself is not be part of the bouncing node's (possibly singleton) group. The result is that the "PTC" message is ignored.

## 6.2 - Quick Communication Interruption Problem

### Example 1 (in 3.2)

In this example, the sender of the "PTC" message has G2 (Group 2) indicated as the last\_stable\_group. Since the destination ("Node 1") does not belong to G2, the "in\_previous\_stable\_group" indication in the "PTC" message is set to "FALSE." Upon receiving the message, Node 1 first sees the value of in\_previous\_stable\_group: FALSE. It then examines whether the sender belongs to last\_stable\_group. For Node 1, last\_stable\_group is G1 (Group 1), and the sender does belong to it (when Node 1 dissolves the group, it forms a singleton group, but it is an unstable one.) Node 1 therefore sees that the two pieces of information are inconsistent. Therefore, Node 1 rejects the message. Node 1 keeps rejecting the "PTC" message until the stability timer expires (typically after about 10 seconds) and Node 1 becomes stable. At this point, Node 1 produces a new last\_stable\_group indication which does not contain the sender of the "PTC." Consequently, the next "PTC" is accepted, since the two pieces of information are consistent. When Node 1 forms a stable singleton group, it sends a node notification saying that all of the other nodes disappeared. And that is the goal: the notification is symmetric to that seen in the other nodes.

### Example 2

In this example, Node 1 has G1 (Group 1) designated as the last\_stable\_group. The other nodes all have G2 (Group 2) designated as the last\_stable\_group. All of the "PTC" messages

have TRUE as an indicator for being "in\_previous\_stable\_group," since all of the recipients belonged to G1. On the other hand, the sender of "PTC" (Node 1) does not belong to G2 (the recipients' last\_stable\_group), so again there is an inconsistency, and the "PTC" is rejected. The same thing happens again until Node 1 forms a singleton stable group.

5

### *Group merge*

Suppose AMG Group 1 (G1) has nodes 1 and 2, while Group 2 (G2) has Nodes 3 and 4. Node 2, which is assumed to be G1's Group Leader sends a "JOIN" message to Node 4, which is G2's Group Leader. Node 4 then sends a "PTC" message to Nodes 1, 2, and 3. For Node 1, the in\_previous\_stable\_group indicator is "FALSE," since Node 1 does not belong to G2. Node 1 itself has Node 4 as not part of the "last\_stable\_group" (G1). The same is true for Node 2. For Node 3, the in\_previous\_stable\_group indicator is "TRUE," since Node 3 belongs to G2. Node 3 itself has Node 4 as part of the "last\_stable\_group" (G2). The end result is that all nodes accept the "PTC", as expected.

### *Real Group Dissolve*

If the Group Leader and Crown Prince fail at the same time, the "third in line" dissolves the group, and all of the adapters in the group form an unstable singleton group. Slowly the remaining members coalesce into a single group. Since the last\_stable\_group indicators contain the group prior to the dissolve, the "PTC" issued during the coalesce phase are accepted.

### *Daemon is blocked*

20 This example actually also falls under Example 2 above. If the daemon is blocked for too long and the adapter was expelled from its AMGs, then the node with the blocked daemon eventually forms a singleton stable group for all of its adapters.

## 7.0 - Discussion in Terms of Related Drawings

The environment in which the present invention is employed is shown in Figure 1. In particular, there is shown a plurality of nodes 100 connected in a network via network adapters 200. Though not specifically shown in the figures herein, the network typically includes a switch which routes messages from source nodes to destination nodes in the network of nodes. Figure 1 also particularly illustrates the possibility that one of the nodes may experience a failure. This failure could occur with the node itself or within the network adapter through which it is coupled to the other nodes.

Figures 2A through 2G illustrate the process in which nodes are added to a node group. This process is initiated by Group Leader #1 (GL1) in a first group sending out a PROCLAIM message to other nodes in a global collection of hardware nodes, as shown in Figure 2A. In particular, the PROCLAIM message is sent to Group Leader #2 in a second, already existing, group of nodes. In general, it is the object to have all of the nodes that count themselves as being in a certain group fully cognizant of all of the other nodes in the group. It is also important to keep in mind that a node may belong to more than one network. The join protocol calls for Group Leader #2 to reply to the PROCLAIM message with a JOIN message which is communicated back to Group Leader #1, as shown in Figure 2B. In accordance with the joining protocol, Group Leader #1 responds to receipt of the JOIN message by sending out a PREPARE\_TO\_COMMIT message (PTC) to all of the nodes that have been identified as part of either of the two groups, as shown in Figure 2C. This includes the fact that Group Leader #1 also effectively sends the PTC message to itself. All of the nodes, including Group Leader #1, responds to the PTC message by sending Group Leader #1 an acknowledgment message (PTC\_ACK). As shown in Figure 2D, this also includes Group Leader #1, as above. In a preferred embodiment of the join protocol, Group Leader #1 sends a "commit to broadcast" message (COMMIT\_BCAST) to a "Mayor" node whose job it is to send an acknowledgment of the COMMIT\_BCAST message to all of the other nodes. Thus, Figure 2F depicts the transmittal of the COMMIT\_BCAST\_ACK message to all of the other nodes, including Group Leader #1. At this point in time, every node in the group is aware of the existence and network address of

every other node that is part of the new node group. Knowledge of the network address is in particular useful in establishing an ordered, closed loop linkage between the nodes in the group. Since the network address is unique and capable of having a numerical ordering from low to high or from high to low, it is an extremely convenient marker for establishing the next neighbor in  
5 the node group for transmission of a heartbeat message. Such a closed loop is illustrated in Figure 2G. The direction of heartbeat message transmission is not critical. It is only important that all of the nodes in the group receive it from an "upstream" node. Although heartbeat messages are sent in a "circular fashion," they are not actually passed along from one node to another in a bucket brigade manner. The timing for receiving a heartbeat message packet and for  
10 sending a heartbeat packet are independent activities. that is to say, a node does not wait to receive a heartbeat message before sending one.

Figure 3A is a simplification of Figure 2G showing the assumed direction for the flow of the heartbeat messages. If the node in the lower right hand portion of Figure 3A fails, as is suggested by the large "X" in Figure 3B, the downstream node, having expected to receive a heartbeat message within a prescribed period of time, informs the Group Leader that his upstream node has died. This is accomplished by sending a "DEATH" message to the Group Leader (GL). The Group Leader now "realizes" that a node has died and that a new node group should be formed.  
15 Upon receipt of the "DEATH" message, the Group Leader sends out a "PREPARE\_TO\_COMMIT" message in the same fashion as shown above in Figure 2C. See  
20 Figure 3C.

Figures 4A, 4B and 4C illustrate the fact that a node may in fact have one or more attached network adapters and thus be connected to more than one node group or node network. It also illustrates the fact that to reach another node it might be the case that one or more network  
25 "hops" might have to be employed. Thus, it becomes very useful for nodes that are connected through a set of network adapters be "aware of" one another. This awareness is accomplished by transmitting a Node Connectivity Message (NCM) to a Group Leader through the Group Leader's adapter. For example, as shown in Figure 4A, connectivity between Node #1 and Node #5 requires a path through two different networks. As shown in Figure 4B, a "Group Connectivity

Message" (GCM) is sent from each Group Leader to all group members. The GCM contains the Adapter Membership Group (AMG) id and the list of nodes that belong to the AMG. Also, for each of these nodes, a list of all "disabled AMG ids" (in the other networks) is included. Figure 4C illustrates the forwarding of the Group Connectivity Message in Network #2.

5 As indicated above, node reachability is maintained across the network by means of a Node  
Connectivity Table. Such a table, with its concomitant entries, is shown in Figure 5A for the  
network and node interconnections shown in Figure 4. In particular, the NCT shown shows two  
groups and two networks. Adapter Membership Group A\_1 includes Node #1, Node #2, Node  
#3 and Node #4. In this group, Node #4 is the Group Leader (GL1). Adapter Membership  
10 Group B\_1 includes Node #3, Node #4, Node #5 and Node #6 with Node #6 being the Group  
Leader (GL2) in AMG B\_1. This structure is reflected in the Node Connectivity Table that  
exists at Node #5, as shown. In particular, the initial Node Connectivity Table shown if Figure  
5A is included to illustrate the state of the system just prior to a node death. In this death  
scenario, Node #2 is assumed to have failed. As shown in Figure 5B, the death of Node #2  
15 results in the formation of a new group from the remaining nodes. Thus Adapter Membership  
Group A\_2 is formed from Node #1, Node #3 and Node #4. The group forming process  
preferably employs the protocols set forth above with particular reference to Figures 2 and 3.  
This group formation is communicated to the other nodes in the network via the Group  
20 Connectivity Message (GCM), as illustrated in Figure 5C. As each node receives the GCM, it  
updates its own Node Connectivity Table. A NCT with updated entries for Node #5 is shown in  
Figure 5C. In this manner, all of the nodes in the network are made aware of group membership  
and "liveness" status.

25 Figure 6 illustrates a scenario (sequence of specific events) in which an inconsistency could  
arise as the result of failure and rapid restart of Node #1. In this scenario it is assumed that the  
time that it takes Network A to detect a node failure, such as through the failure of the heartbeat  
daemon to receive a heartbeat message from its upstream neighbor, is less than the corresponding  
time for this same operation in Network B. For example, such a situation could easily arise if the  
networks are set up with different tunable parameters, such as the default adapter failure

detection time-out value. In the example shown, it is assumed that the failure detection time for Network A is 10 seconds while the corresponding failure detection time for Network B is 40 seconds. In this scenario it is also assumed that Node #1 and Node #2 are part of an Adapter Membership Group on Networks A and B.

5        In Figure 6 it is seen that at time T = 0 (seconds being the assumed time unit), termination of the heartbeat daemon on Node #1 results in missing heartbeats occurring for this node in Networks A and B. However, because of the slower detection process on Network B, the death of Node #1 is not recognized until T = 40. Nonetheless, on Network A, Node #1 is detected as having "died" at T = 10. At T = 20, the heartbeat daemon is restarted on Node #1 (this is the start of the "quick restart" phenomenon that gives rise to one of the problems considered herein; it should also be noted that the use of the adjective "quick" to describe the restart is meant to be a relative term: it is quick relative to the timing for node and/or adapter failure detection in a connected network of nodes.) With the "rebirth" of Node #1 in Network A it then joins a group on Network A at T = 25. At T = 30, Node #1 is still seen as being alive on Network B which has not as yet determined that it has failed. At T = 40, the death of Node #1 is finally detected at Node #2. Node #1 is now "free" to join a group on Network B which Node #2 is capable of processing at T = 45. At this point we potentially have Node #1 as part of a group in Network A and also a member of a different group in Network B. Nowhere in the sequence from T = 0 to T = 45 shown in Figure 6 is Node #1 seen as being completely unreachable in both networks.  
10  
15  
20  
25  
30  
35  
40  
45  
50  
55  
60  
65  
70  
75  
80  
85  
90  
95  
100  
105  
110  
115  
120  
125  
130  
135  
140  
145  
150  
155  
160  
165  
170  
175  
180  
185  
190  
195  
200  
205  
210  
215  
220  
225  
230  
235  
240  
245  
250  
255  
260  
265  
270  
275  
280  
285  
290  
295  
300  
305  
310  
315  
320  
325  
330  
335  
340  
345  
350  
355  
360  
365  
370  
375  
380  
385  
390  
395  
400  
405  
410  
415  
420  
425  
430  
435  
440  
445  
450  
455  
460  
465  
470  
475  
480  
485  
490  
495  
500  
505  
510  
515  
520  
525  
530  
535  
540  
545  
550  
555  
560  
565  
570  
575  
580  
585  
590  
595  
600  
605  
610  
615  
620  
625  
630  
635  
640  
645  
650  
655  
660  
665  
670  
675  
680  
685  
690  
695  
700  
705  
710  
715  
720  
725  
730  
735  
740  
745  
750  
755  
760  
765  
770  
775  
780  
785  
790  
795  
800  
805  
810  
815  
820  
825  
830  
835  
840  
845  
850  
855  
860  
865  
870  
875  
880  
885  
890  
895  
900  
905  
910  
915  
920  
925  
930  
935  
940  
945  
950  
955  
960  
965  
970  
975  
980  
985  
990  
995  
1000  
1005  
1010  
1015  
1020  
1025  
1030  
1035  
1040  
1045  
1050  
1055  
1060  
1065  
1070  
1075  
1080  
1085  
1090  
1095  
1100  
1105  
1110  
1115  
1120  
1125  
1130  
1135  
1140  
1145  
1150  
1155  
1160  
1165  
1170  
1175  
1180  
1185  
1190  
1195  
1200  
1205  
1210  
1215  
1220  
1225  
1230  
1235  
1240  
1245  
1250  
1255  
1260  
1265  
1270  
1275  
1280  
1285  
1290  
1295  
1300  
1305  
1310  
1315  
1320  
1325  
1330  
1335  
1340  
1345  
1350  
1355  
1360  
1365  
1370  
1375  
1380  
1385  
1390  
1395  
1400  
1405  
1410  
1415  
1420  
1425  
1430  
1435  
1440  
1445  
1450  
1455  
1460  
1465  
1470  
1475  
1480  
1485  
1490  
1495  
1500  
1505  
1510  
1515  
1520  
1525  
1530  
1535  
1540  
1545  
1550  
1555  
1560  
1565  
1570  
1575  
1580  
1585  
1590  
1595  
1600  
1605  
1610  
1615  
1620  
1625  
1630  
1635  
1640  
1645  
1650  
1655  
1660  
1665  
1670  
1675  
1680  
1685  
1690  
1695  
1700  
1705  
1710  
1715  
1720  
1725  
1730  
1735  
1740  
1745  
1750  
1755  
1760  
1765  
1770  
1775  
1780  
1785  
1790  
1795  
1800  
1805  
1810  
1815  
1820  
1825  
1830  
1835  
1840  
1845  
1850  
1855  
1860  
1865  
1870  
1875  
1880  
1885  
1890  
1895  
1900  
1905  
1910  
1915  
1920  
1925  
1930  
1935  
1940  
1945  
1950  
1955  
1960  
1965  
1970  
1975  
1980  
1985  
1990  
1995  
2000  
2005  
2010  
2015  
2020  
2025  
2030  
2035  
2040  
2045  
2050  
2055  
2060  
2065  
2070  
2075  
2080  
2085  
2090  
2095  
2100  
2105  
2110  
2115  
2120  
2125  
2130  
2135  
2140  
2145  
2150  
2155  
2160  
2165  
2170  
2175  
2180  
2185  
2190  
2195  
2200  
2205  
2210  
2215  
2220  
2225  
2230  
2235  
2240  
2245  
2250  
2255  
2260  
2265  
2270  
2275  
2280  
2285  
2290  
2295  
2300  
2305  
2310  
2315  
2320  
2325  
2330  
2335  
2340  
2345  
2350  
2355  
2360  
2365  
2370  
2375  
2380  
2385  
2390  
2395  
2400  
2405  
2410  
2415  
2420  
2425  
2430  
2435  
2440  
2445  
2450  
2455  
2460  
2465  
2470  
2475  
2480  
2485  
2490  
2495  
2500  
2505  
2510  
2515  
2520  
2525  
2530  
2535  
2540  
2545  
2550  
2555  
2560  
2565  
2570  
2575  
2580  
2585  
2590  
2595  
2600  
2605  
2610  
2615  
2620  
2625  
2630  
2635  
2640  
2645  
2650  
2655  
2660  
2665  
2670  
2675  
2680  
2685  
2690  
2695  
2700  
2705  
2710  
2715  
2720  
2725  
2730  
2735  
2740  
2745  
2750  
2755  
2760  
2765  
2770  
2775  
2780  
2785  
2790  
2795  
2800  
2805  
2810  
2815  
2820  
2825  
2830  
2835  
2840  
2845  
2850  
2855  
2860  
2865  
2870  
2875  
2880  
2885  
2890  
2895  
2900  
2905  
2910  
2915  
2920  
2925  
2930  
2935  
2940  
2945  
2950  
2955  
2960  
2965  
2970  
2975  
2980  
2985  
2990  
2995  
3000  
3005  
3010  
3015  
3020  
3025  
3030  
3035  
3040  
3045  
3050  
3055  
3060  
3065  
3070  
3075  
3080  
3085  
3090  
3095  
3100  
3105  
3110  
3115  
3120  
3125  
3130  
3135  
3140  
3145  
3150  
3155  
3160  
3165  
3170  
3175  
3180  
3185  
3190  
3195  
3200  
3205  
3210  
3215  
3220  
3225  
3230  
3235  
3240  
3245  
3250  
3255  
3260  
3265  
3270  
3275  
3280  
3285  
3290  
3295  
3300  
3305  
3310  
3315  
3320  
3325  
3330  
3335  
3340  
3345  
3350  
3355  
3360  
3365  
3370  
3375  
3380  
3385  
3390  
3395  
3400  
3405  
3410  
3415  
3420  
3425  
3430  
3435  
3440  
3445  
3450  
3455  
3460  
3465  
3470  
3475  
3480  
3485  
3490  
3495  
3500  
3505  
3510  
3515  
3520  
3525  
3530  
3535  
3540  
3545  
3550  
3555  
3560  
3565  
3570  
3575  
3580  
3585  
3590  
3595  
3600  
3605  
3610  
3615  
3620  
3625  
3630  
3635  
3640  
3645  
3650  
3655  
3660  
3665  
3670  
3675  
3680  
3685  
3690  
3695  
3700  
3705  
3710  
3715  
3720  
3725  
3730  
3735  
3740  
3745  
3750  
3755  
3760  
3765  
3770  
3775  
3780  
3785  
3790  
3795  
3800  
3805  
3810  
3815  
3820  
3825  
3830  
3835  
3840  
3845  
3850  
3855  
3860  
3865  
3870  
3875  
3880  
3885  
3890  
3895  
3900  
3905  
3910  
3915  
3920  
3925  
3930  
3935  
3940  
3945  
3950  
3955  
3960  
3965  
3970  
3975  
3980  
3985  
3990  
3995  
4000  
4005  
4010  
4015  
4020  
4025  
4030  
4035  
4040  
4045  
4050  
4055  
4060  
4065  
4070  
4075  
4080  
4085  
4090  
4095  
4100  
4105  
4110  
4115  
4120  
4125  
4130  
4135  
4140  
4145  
4150  
4155  
4160  
4165  
4170  
4175  
4180  
4185  
4190  
4195  
4200  
4205  
4210  
4215  
4220  
4225  
4230  
4235  
4240  
4245  
4250  
4255  
4260  
4265  
4270  
4275  
4280  
4285  
4290  
4295  
4300  
4305  
4310  
4315  
4320  
4325  
4330  
4335  
4340  
4345  
4350  
4355  
4360  
4365  
4370  
4375  
4380  
4385  
4390  
4395  
4400  
4405  
4410  
4415  
4420  
4425  
4430  
4435  
4440  
4445  
4450  
4455  
4460  
4465  
4470  
4475  
4480  
4485  
4490  
4495  
4500  
4505  
4510  
4515  
4520  
4525  
4530  
4535  
4540  
4545  
4550  
4555  
4560  
4565  
4570  
4575  
4580  
4585  
4590  
4595  
4600  
4605  
4610  
4615  
4620  
4625  
4630  
4635  
4640  
4645  
4650  
4655  
4660  
4665  
4670  
4675  
4680  
4685  
4690  
4695  
4700  
4705  
4710  
4715  
4720  
4725  
4730  
4735  
4740  
4745  
4750  
4755  
4760  
4765  
4770  
4775  
4780  
4785  
4790  
4795  
4800  
4805  
4810  
4815  
4820  
4825  
4830  
4835  
4840  
4845  
4850  
4855  
4860  
4865  
4870  
4875  
4880  
4885  
4890  
4895  
4900  
4905  
4910  
4915  
4920  
4925  
4930  
4935  
4940  
4945  
4950  
4955  
4960  
4965  
4970  
4975  
4980  
4985  
4990  
4995  
5000  
5005  
5010  
5015  
5020  
5025  
5030  
5035  
5040  
5045  
5050  
5055  
5060  
5065  
5070  
5075  
5080  
5085  
5090  
5095  
5100  
5105  
5110  
5115  
5120  
5125  
5130  
5135  
5140  
5145  
5150  
5155  
5160  
5165  
5170  
5175  
5180  
5185  
5190  
5195  
5200  
5205  
5210  
5215  
5220  
5225  
5230  
5235  
5240  
5245  
5250  
5255  
5260  
5265  
5270  
5275  
5280  
5285  
5290  
5295  
5300  
5305  
5310  
5315  
5320  
5325  
5330  
5335  
5340  
5345  
5350  
5355  
5360  
5365  
5370  
5375  
5380  
5385  
5390  
5395  
5400  
5405  
5410  
5415  
5420  
5425  
5430  
5435  
5440  
5445  
5450  
5455  
5460  
5465  
5470  
5475  
5480  
5485  
5490  
5495  
5500  
5505  
5510  
5515  
5520  
5525  
5530  
5535  
5540  
5545  
5550  
5555  
5560  
5565  
5570  
5575  
5580  
5585  
5590  
5595  
5600  
5605  
5610  
5615  
5620  
5625  
5630  
5635  
5640  
5645  
5650  
5655  
5660  
5665  
5670  
5675  
5680  
5685  
5690  
5695  
5700  
5705  
5710  
5715  
5720  
5725  
5730  
5735  
5740  
5745  
5750  
5755  
5760  
5765  
5770  
5775  
5780  
5785  
5790  
5795  
5800  
5805  
5810  
5815  
5820  
5825  
5830  
5835  
5840  
5845  
5850  
5855  
5860  
5865  
5870  
5875  
5880  
5885  
5890  
5895  
5900  
5905  
5910  
5915  
5920  
5925  
5930  
5935  
5940  
5945  
5950  
5955  
5960  
5965  
5970  
5975  
5980  
5985  
5990  
5995  
6000  
6005  
6010  
6015  
6020  
6025  
6030  
6035  
6040  
6045  
6050  
6055  
6060  
6065  
6070  
6075  
6080  
6085  
6090  
6095  
6100  
6105  
6110  
6115  
6120  
6125  
6130  
6135  
6140  
6145  
6150  
6155  
6160  
6165  
6170  
6175  
6180  
6185  
6190  
6195  
6200  
6205  
6210  
6215  
6220  
6225  
6230  
6235  
6240  
6245  
6250  
6255  
6260  
6265  
6270  
6275  
6280  
6285  
6290  
6295  
6300  
6305  
6310  
6315  
6320  
6325  
6330  
6335  
6340  
6345  
6350  
6355  
6360  
6365  
6370  
6375  
6380  
6385  
6390  
6395  
6400  
6405  
6410  
6415  
6420  
6425  
6430  
6435  
6440  
6445  
6450  
6455  
6460  
6465  
6470  
6475  
6480  
6485  
6490  
6495  
6500  
6505  
6510  
6515  
6520  
6525  
6530  
6535  
6540  
6545  
6550  
6555  
6560  
6565  
6570  
6575  
6580  
6585  
6590  
6595  
6600  
6605  
6610  
6615  
6620  
6625  
6630  
6635  
6640  
6645  
6650  
6655  
6660  
6665  
6670  
6675  
6680  
6685  
6690  
6695  
6700  
6705  
6710  
6715  
6720  
6725  
6730  
6735  
6740  
6745  
6750  
6755  
6760  
6765  
6770  
6775  
6780  
6785  
6790  
6795  
6800  
6805  
6810  
6815  
6820  
6825  
6830  
6835  
6840  
6845  
6850  
6855  
6860  
6865  
6870  
6875  
6880  
6885  
6890  
6895  
6900  
6905  
6910  
6915  
6920  
6925  
6930  
6935  
6940  
6945  
6950  
6955  
6960  
6965  
6970  
6975  
6980  
6985  
6990  
6995  
7000  
7005  
7010  
7015  
7020  
7025  
7030  
7035  
7040  
7045  
7050  
7055  
7060  
7065  
7070  
7075  
7080  
7085  
7090  
7095  
7100  
7105  
7110  
7115  
7120  
7125  
7130  
7135  
7140  
7145  
7150  
7155  
7160  
7165  
7170  
7175  
7180  
7185  
7190  
7195  
7200  
7205  
7210  
7215  
7220  
7225  
7230  
7235  
7240  
7245  
7250  
7255  
7260  
7265  
7270  
7275  
7280  
7285  
7290  
7295  
7300  
7305  
7310  
7315  
7320  
7325  
7330  
7335  
7340  
7345  
7350  
7355  
7360  
7365  
7370  
7375  
7380  
7385  
7390  
7395  
7400  
7405  
7410  
7415  
7420  
7425  
7430  
7435  
7440  
7445  
7450  
7455  
7460  
7465  
7470  
7475  
7480  
7485  
7490  
7495  
7500  
7505  
7510  
7515  
7520  
7525  
7530  
7535  
7540  
7545  
7550  
7555  
7560  
7565  
7570  
7575  
7580  
7585  
7590  
7595  
7600  
7605  
7610  
7615  
7620  
7625  
7630  
7635  
7640  
7645  
7650  
7655  
7660  
7665  
7670  
7675  
7680  
7685  
7690  
7695  
7700  
7705  
7710  
7715  
7720  
7725  
7730  
7735  
7740  
7745  
7750  
7755  
7760  
7765  
7770  
7775  
7780  
7785  
7790  
7795  
7800  
7805  
7810  
7815  
7820  
7825  
7830  
7835  
7840  
7845  
7850  
7855  
7860  
7865  
7870  
7875  
7880  
7885  
7890  
7895  
7900  
7905  
7910  
7915  
7920  
7925  
7930  
7935  
7940  
7945  
7950  
7955  
7960  
7965  
7970  
7975  
7980  
7985  
7990  
7995  
8000  
8005  
8010  
8015  
8020  
8025  
8030  
8035  
8040  
8045  
8050  
8055  
8060  
8065  
8070  
8075  
8080  
8085  
8090  
8095  
8100  
8105  
8110  
8115  
8120  
8125  
8130  
8135  
8140  
8145  
8150  
8155  
8160  
8165  
8170  
8175  
8180  
8185  
8190  
8195  
8200  
8205  
8210  
8215  
8220  
8225  
8230  
8235  
8240  
8245  
8250  
8255  
8260  
8265  
8270  
8275  
8280  
8285  
8290  
8295  
8300  
8305  
8310  
8315  
8320  
8325  
8330  
8335  
8340  
8345  
8350  
8355  
8360  
8365  
8370  
8375  
8380  
8385  
8390  
8395  
8400  
8405  
8410  
8415  
8420  
8425  
8430  
8435  
8440  
8445  
8450  
8455  
8460  
8465  
8470  
8475  
8480  
8485  
8490  
8495  
8500  
8505  
8510  
8515  
8520  
8525  
8530  
8535  
8540  
8545  
8550  
8555  
8560  
8565  
8570  
8575  
8580  
8585  
8590  
8595  
8600  
8605  
8610  
8615  
8620  
8625  
8630  
8635  
8640  
8645  
8650  
8655  
8660  
8665  
8670  
8675  
8680  
8685  
8690  
8695  
8700  
8705  
8710  
8715  
8720  
8725  
8730  
8735  
8740  
8745  
8750  
8755  
8760  
8765  
8770  
8775  
8780  
8785  
8790  
8795  
8800  
8805  
8810  
8815  
8820  
8825  
8830  
8835  
8840  
8845  
8850  
8855  
8860  
8865  
8870  
8875  
8880  
8885  
8890  
8895  
8900  
8905  
8910  
8915  
8920  
8925  
8930  
8935  
8940  
8945  
8950  
8955  
8960  
8965  
8970  
8975  
8980  
8985  
8990  
8995  
9000  
9005  
9010  
9015  
9020  
9025  
9030  
9035  
9040  
9045  
9050  
9055  
9060  
9065  
9070  
9075  
9080  
9085  
9090  
9095  
9100  
9105  
9110  
9115  
9120  
9125  
9130  
9135  
9140  
9145  
9150  
9155  
9160  
9165  
9170  
9175  
9180  
9185  
9190  
9195  
9200  
9205  
9210  
9215  
9220  
9225  
9230  
9235  
9240  
9245  
9250  
9255  
9260  
9265  
9270  
9

experienced after Node #1 is expelled from the group. At Node #1, the lack of a heartbeat message and responses from either the Group Leader or from the Crown Prince nodes forces Node #1 into the formation of a singleton group (a group with only one member). Such groups are inherently "unstable" as that term is used herein where it refers to the likelihood that any  
5 given node is soon apt to change its membership status. If and when the adapter problems at Node #1 end, it resumes normal operations and is able to rejoin the group. Because Node #1 formed a singleton unstable group, Node #1 does not perceive any of the other nodes as unreachable. At the same time, the other nodes do see Node #1 as reachable, thus producing an undesirable state of reachability inconsistency.

Figure 8 is similar to Figure 7 except that it illustrates the sequence of events that would occur in Figure 7 if Node #1 had been either the Group Leader or the Crown Prince Node. As in Figure 7, it is here assumed that the network adapter at Node #1 experiences a failure in which it is able to receive message packets but is not able to send them. When Node #2 and the other nodes in its group detect the fact that heartbeat messages are not being sent from the Group  
10 Leader (Node #1), Node #2 and the other nodes form a new Adapter Membership Group, without Node #1. As a result, no heartbeat messages are sent to Node #1 since it is not part of the new group formed with Node #2. Therefore, Node #1 no longer receives heartbeat messages from its upstream neighbor, which it assumes has died. Node #1, still "thinking" that it has Group Leader status sends a PREPARE\_TO\_COMMIT (PTC) message to all of its network peers, except for  
15 the node which it assumes is its dead upstream neighbor. The group containing Node #2 replies with a PTC\_ACK message and Node #1 replies to this message by sending a COMMIT message to all of its network peer nodes. Node #1 is thereafter made part of a new group (that is, the group formed by the protocol, as described above which, from Node #1's point of view is actually a DEATH protocol as opposed to a JOIN protocol. However, in this example, Node #1  
20 never sees the other nodes as being unreachable while the other nodes do see Node #1 as being unreachable for the period of time illustrated (Node #1 dead).

Figure 9 depicts how Adapter IDs and Group IDs are created. Adapter IDs are formed with a pair of indicators: (IP address; instance number). The IP address is the adapter's IP address. The

instance number is assigned when the adapter is initialized or reinitialized into a singleton group. This number is any uniquely selected number increasing with time; in particular, it is preferably derived from the time-of-day value. As also seen in Figure 9, the Group ID also comprises a pair of indicators: (IP address; instance number). The address refers to the IP address of the adapter which is the Group Leader. The instance number is assigned at the time the group is formed.  
5 Each member of an AMG has a distinct Adapter ID, while all members of the group have the same Group ID.

Figure 10 illustrates a preferable format for the packet that is sent from one daemon to another over the network. The packet includes a message type ("HEART BEAT," "PROCLAIM," etc.), the Adapter ID of the source of the message, the Group ID of the source of the message, the Adapter and Group ID of the destination of the message, and finally a "payload" area which is specific to each message type.

Figure 11 illustrates how Adapter and Group IDs change over time in a scenario where the daemon at Node #1 terminates and is later restarted. When the daemon at Node #1 terminates, initially there is no change in Node #2's view, since not enough heartbeats are missed from Node #1. When the daemon at Node #1 is restarted, it has a new Adapter ID, Group ID, and the local adapter forms a singleton AMG.

Figure 12 builds on the scenario illustrated in Figure 11. The adapter on Node #1 sends a "PROCLAIM" message that eventually reaches Node #2. The message reveals the following:

- 20
- while the address part of the "source Adapter ID" ("1.1.1.1") can actually be found in Node #2's AMG, the instance number portions ("7259" in the AMG and "7901" in the message) do not match; and
  - the Source Group ID in the message does not match the Group ID stored at Node #2.

The indication provided by the inconsistencies above is enough for Node #2 to determine that Node #1 bounced.

Figures 13A and 13B illustrate another mechanism used to detect the occurrence of a bounce. Initially (see Figure 13A), the daemon at one of the nodes (that with the large "X") terminates, and therefore no more "HEART BEAT" messages are sent by it. The daemon is restarted on the same node. The daemon then starts receiving "HEART BEAT" messages from the previous instance's upstream neighbor (see Figure 13B). Upon receiving such messages, the node notices that:

- the messages are not intended for the node's current Group ID; and
- the current group is singleton.

The inconsistencies above are interpreted as the recipient node having bounced. In response to the "HEART BEAT" message, it sends a "NOT YOUR NEIGHBOR" response back to the sender, which in turn sends a "DEATH" message to the Group Leader, informing it about the demise of the bounced node.

Figures 14A through 14D illustrate an example of the Quick Communication Interruption problem. Initially (see Figure 14A), three nodes belong to the AMG. Note that the Group IDs, AMG, and "last\_stable\_group" are the same across all three nodes. At this point, the adapter at Node #3 (the Group Leader) is assumed to suffer a temporary outage. Figure 14B shows the results of Node #2 (the new Group Leader) expelling Node #3 from the group with the result that the Group ID, AMG, and "last\_stable\_group" all change to reflect the formation of the new group. Meanwhile, Node #3 has still not missed enough "HEART BEAT" messages to take any action. In Figure 14C, Node #3 finally misses enough "HEART BEAT" messages to declare its upstream neighbor "down" and to send a "PTC" message to Node #2. The PTC message has its "in\_previous\_stable\_group" indicator set equal to 1, since the recipient ("1.1.1.2") is still listed as part of "last\_stable\_group" on Node #3. On the other hand, the sender of the message

("1.1.1.3") does not belong to the recipient's "last\_stable\_group." The inconsistency between the latter and the value "1" in the value of "in\_previous\_stable\_group" leads the recipient to reject the message. Since all of the PTCs sent by Node #3 are rejected, Node #3 eventually forms a singleton group (see Figure 14D), at which points it declares all remaining nodes as down. The  
5 next PTC sent by Node #3 has its "in\_previous\_stable\_group" indicator set equal to 0, which prompts Node #2 to accept the message. Thus, the goal of having Node #3 perceive all other nodes as down (in symmetry to Node #3 itself having been seen as down by the others) is achieved.

Figures 15A through 15E are similar to Figure 14, but they focus on Node #1, which is not the Group Leader. Initially (see Figure 15A) all the nodes have the same view of the group membership. When an outage occurs on Node #1's adapter, it is expelled from Node #3's group, as shown in Figure 15B. In this figure, the instance number of Node #3's Group ID is changed to reflect the new group which excludes Node #1. Figure 15C shows Node #1 dissolving its AMG - - after it started missing "HEART BEAT" messages, sent a "DEATH" message, but did not get a "PTC" message. While Node #1 forms a new AMG with a new Group ID, the AMG is unstable and therefore does not cause a change in "last\_stable\_group." When Node #1's adapter recovers, a "PTC" message from Node #3 finally reaches Node #1 (Figure 15D). The  
10 "in\_previous\_stable\_group" indicator in the message has value 0, since Node #1 does not belong to Node #3's "last\_stable\_group." On the other hand, Node #1 still has Node #3 as belonging to its version of "last\_stable\_group." Therefore again there is an inconsistency, which prompts  
20 Node #1 to reject the "PTC" message. Eventually (see Figure 15E), Node #1 forms a stable AMG, which finally causes "last\_stable\_group" to change. Because then the sender of the "PTC" is no longer part of Node #1's "last\_stable\_group," the inconsistency disappears, and the "PTC" is finally accepted. Thus, here too, the goal of having Node #1 see all of the other nodes as  
25 "down" (since all the other nodes saw Node #1 as down) is now satisfied.

While the invention has been described in detail herein in accordance with certain preferred embodiments thereof, many modifications and changes therein may be effected by

those skilled in the art. Accordingly, it is intended by the appended claims to cover all such modifications and changes as fall within the true spirit and scope of the invention

POU920020016US1